

Topic 5

 Unconstrained minimization.

$$D. \quad \min f(x)$$

Assume $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, twice differentiable

Assume that the problem is solvable, that is, an x^* exists.

Denote

$$p^* = f(x^*) = \inf_{x \in \mathbb{R}^n} f(x)$$

Since $f(x)$ is differentiable and convex, for x^* to be optimal it is necessary and sufficient that

$$\nabla f(x^*) = 0.$$

If this equation has analytical solution, we are done!

Often the above equation does not allow an analytical solution and hence must be solved numerically, through an iterative algorithm.

This is an algorithm that computes a sequence of points

$$x^{(0)}, x^{(1)}, x^{(2)}, \dots \in \text{dom } f$$

such that $f(x^{(k)}) \rightarrow p^*$ as $k \rightarrow \infty$.

Such an sequence is called a minimizing sequence for the original optimization problem.

The algorithm terminates when $f(x^{(k)}) - p^* \leq \epsilon$ where $\epsilon > 0$ is some pre-specified tolerance.

+) Initial point and sublevel set:

The methods we are going to study for unconstrained minimization, here require a suitable starting point $x^{(0)}$ such that:

$$\left\{ \begin{array}{l} x^{(0)} \in \text{dom } f \\ \text{sublevel set } S = \{x \in \text{dom } f \mid f(x) \leq f(x^{(0)})\} \text{ is closed.} \end{array} \right.$$

Examples: (i) Least-square

$$\min \|Ax - b\|_2^2 = x^T(A^T A)x - 2(A^T b)^T x + b^T b$$

The optimality condition is

$$A^T A x^* = A^T b.$$

- If $A^T A$ is invertible, then the problem has a unique solution.
- If $A^T A \succeq 0$ but is not strictly $\succ 0$, then any solution of the above equation is optimal.
- If the optimality equation does not have a solution then the optimization problem is unbounded below ($p^* = -\infty$).

(ii) GP: $\min f(x) = \log\left(\sum_{i=1}^m e^{a_i^T x + b_i}\right)$

The optimality condition is

$$\nabla f(x^*) = \frac{1}{\sum_{i=1}^m e^{a_i^T x^* + b_i}} \sum_{i=1}^m e^{a_i^T x^* + b_i} \cdot a_i = 0$$

This equation does not have an analytical solution in general so we must rely on an iterative algorithm to find x^* .

*) Notes: The condition on closed sublevel sets is usually hard to verify, but holds if
dom $f = \mathbb{R}^n$ (whole space)
or $f(x) \rightarrow \infty$ as $x \rightarrow$ boundary of dom f .

+) Strong convexity and implications:

Recall that the condition for $f(x)$ to be convex is
 $\nabla^2 f(x) \succeq 0$ (provided $f(x)$ is differentiable).

Here we assume strong convexity: $\exists m > 0$ s.t.

$$\nabla^2 f(x) \succeq m I, \quad \forall x \in S.$$

That is, the Hessian is strictly positive definite (on set S).

By the Taylor's expansion and mean-value theorem:

$$\forall x, y \in S: \exists z \in [x, y] \text{ s.t.}$$

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x)$$

By strong convexity then

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} (y-x)^T (y-x)$$

$$\forall x, y \in S$$

When $m=0$, we recover the first order condition for convexity.

For $m > 0$, we obtain a better bound.

We can obtain a bound on p^* as follows:

Minimize the RHS expression wrt y :

$$\text{minimizer } \tilde{y} = x - \frac{1}{m} \nabla f(x) \quad (\text{setting the derivatives wrt } y \text{ to zero})$$

Hence

$$f(y) \geq f(x) + \nabla f(x)^T (\tilde{y}-x) + \frac{m}{2} (\tilde{y}-x)^T (\tilde{y}-x)$$

$$= f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2, \quad \forall y \in S.$$

$$\rightarrow p^* \geq f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2.$$

This inequality shows that if $\|\nabla f(x)\|_2$ is small at a point then the point is nearly optimal.

If we know m , we can also use it as a stopping criterion.

$$f(x) - p^* \leq \frac{1}{2m} \|\nabla f(x)\|_2^2.$$

We can also show a bound on the optimal point:

Let $y = x^*$

$$\begin{aligned} \rightarrow p^* = f(x^*) &\geq f(x) + \nabla f(x)^T (x^* - x) + \frac{m}{2} \|x^* - x\|_2^2 \\ &\geq f(x) - \|\nabla f(x)\| \|x^* - x\| + \frac{m}{2} \|x^* - x\|_2^2 \end{aligned}$$

(by Cauchy-Schwartz)

but $p^* \leq f(x) \forall x$

$$\rightarrow -\|\nabla f(x)\| \cdot \|x^* - x\|_2 + \frac{m}{2} \|x^* - x\|_2^2 \leq 0 \quad \forall x$$

$$\rightarrow \|x - x^*\|_2 \leq \frac{2}{m} \|\nabla f(x)\|_2 \quad \forall x.$$

One implication of this is the optimal point x^* is unique.

→ Upper bound on $\nabla^2 f(x)$:

Strong convexity implies that the sublevelsets contained in S , thus S is bounded.

Thus $\exists M > 0$ such that

$$\nabla^2 f(x) \preceq MI. \quad \forall x \in S.$$

This implies

$$\forall x, y \in S: f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{M}{2} \|y - x\|_2^2$$

which, by minimizing both sides over y , yields

$$p^* \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2$$

Thus

$$f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2 \leq p^* \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2 \quad \forall x \in S.$$

+) Condition number bound:

$$mI \leq \nabla^2 f(x) \leq MI \quad \forall x \in S.$$

The ratio $K = \frac{M}{m}$ is an upper bound on the condition number of $\nabla^2 f(x)$.

This condition number has a strong effect on the efficiency of some common methods for unconstrained optimization.

o Usually the constants m and M are unknown, so they cannot be used directly as a practical stopping criterion.

However, they can be used as a conceptual stopping criterion and used in convergence proofs for algorithms. These convergence proofs usually assume some (unknown) constants m, M , except for a special class of convex functions (self-concordant).
More on this later.

2). Descent methods.

o Consider algorithms that produce a minimizing sequence $x^{(k)}$

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$$

where $\Delta x^{(k)}$: step or search direction (vector in \mathbb{R}^n)

k : iteration number

$t^{(k)} > 0$: step size or step length ($t^{(k)} = 0$ if $x^{(k)}$ is optimal)

o When focusing on one iteration, we write

$$x^+ = x + t \Delta x.$$

+) Condition number bound:

$$mI \leq \nabla^2 f(x) \leq MI \quad \forall x \in S.$$

The ratio $K = \frac{M}{m}$ is an upper bound on the condition number of $\nabla^2 f(x)$.

This condition number has a strong effect on the efficiency of some common methods for unconstrained optimization.

o Usually the constants m and M are unknown, so they cannot be used directly as a practical stopping criterion.

However, they can be used as a conceptual stopping criterion and used in convergence proofs for algorithms. These convergence proofs usually assume some (unknown) constants m, M , except for a special class of convex functions (self-concordant).
More on this later.

ecture 17:

2). Descent methods.

o Consider algorithms that produce a minimizing sequence $x^{(k)}$:

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$$

where $\Delta x^{(k)}$: step or search direction (vector in \mathbb{R}^n)

k : iteration number

$t^{(k)} > 0$: step size or step length ($t^{(k)} = 0$ if $x^{(k)}$ is optimal)

o When focusing on one iteration, we write

$$x^+ = x + t \Delta x.$$

$$f(x) - p^* \leq \frac{1}{2m} \|\nabla f(x)\|_2^2.$$

We can also show a bound on the optimal point:

Let $y = x^*$

$$\begin{aligned} \rightarrow p^* = f(x^*) &\geq f(x) + \nabla f(x)^T (x^* - x) + \frac{m}{2} \|x^* - x\|_2^2 \\ &\geq f(x) - \|\nabla f(x)\| \|x^* - x\| + \frac{m}{2} \|x^* - x\|_2^2 \\ &\quad \text{(by Cauchy-Schwartz)} \end{aligned}$$

but $p^* \leq f(x) \forall x$

$$\rightarrow -\|\nabla f(x)\| \cdot \|x^* - x\|_2 + \frac{m}{2} \|x^* - x\|_2^2 \leq 0 \quad \forall x$$

$$\rightarrow \|x - x^*\|_2 \leq \frac{2}{m} \|\nabla f(x)\|_2 \quad \forall x.$$

One implication of this is the optimal point x^* is unique.

+) Upper bound on $\nabla^2 f(x)$:

Strong convexity implies that the sublevelsets contained in S , thus S is bounded.

Thus $\exists M > 0$ such that

$$\nabla^2 f(x) \preceq MI. \quad \forall x \in S.$$

This implies

$$\forall x, y \in S: f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{M}{2} \|y - x\|_2^2$$

which, by minimizing both sides over y , yields

$$p^* \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2$$

Thus

$$f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2 \leq p^* \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2 \quad \forall x \in S.$$

A method is called a descent method if
 $f(x^{(k+1)}) < f(x^{(k)})$
 except when $x^{(k)}$ is optimal.

From convexity of $f(x)$, we know that

$$\nabla f(x^{(k)})^T (y - x^{(k)}) \geq 0 \text{ implies } f(y) \geq f(x^{(k)})$$

since

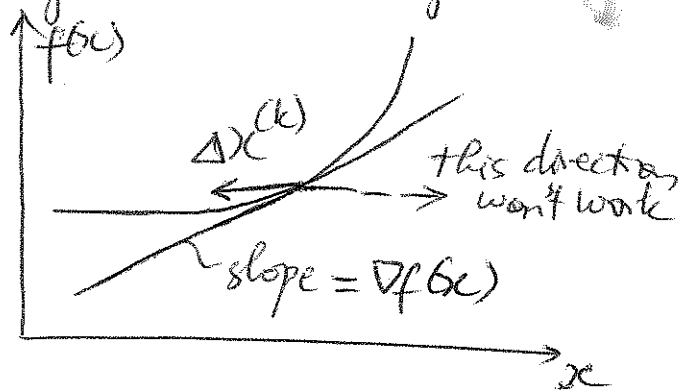
$$f(y) \geq f(x^{(k)}) + \nabla f(x^{(k)})^T (y - x^{(k)}) \quad \forall y, x^{(k)} \in \text{dom} f$$

(1st order condition)

Thus the search direction in a descent method must satisfy

$$\nabla f(x^{(k)})^T \Delta x^{(k)} < 0$$

→ it must make an acute angle with the negative gradient at $x^{(k)}$.



→ General descent method.

given a starting point $x \in \text{dom} f$

repeat 1. Determine a descent direction

2. Line search: choose a step size $t > 0$

3. Update: $x := x + t \Delta x$

until stopping criterion is satisfied.

The stopping criterion is often of the form

$$\|\nabla f(x)\|_2 < \epsilon$$

where ϵ is a small, positive, prechosen tolerance value

+ Line search: Line search is used to select the stepsize t , which determine how far along the line $x + t\Delta x$ the next iterate will be.

o Exact line search:

$$t = \arg \min_{s \geq 0} f(x + s\Delta x)$$

This method can be used if the minimizer of $f(x)$ along the line $x + s\Delta x$ can be found analytically, or if the cost of finding the exact minimum is low.

o Backtracking line search:

Most practical line searches are inexact: to reduce $f(x)$ "enough" along the direction Δx .

Backtracking line search: 2 constant parameters $\alpha \in (0, \frac{1}{2})$, $\beta \in ($

(starting at $t=1$)

$$\text{while } f(x + t\Delta x) > f(x) + \alpha t \nabla f(x)^T \Delta x$$

$$t = \beta t$$

(Typical values: $\alpha \in [0.01, 0.3]$, $\beta \in [0.1, 0.8]$).

o Since Δx is a descent direction, then $\nabla f(x)^T \Delta x < 0$.

Thus for t small enough we have

$$f(x + t\Delta x) \approx f(x) + t \nabla f(x)^T \Delta x$$

$$< f(x) + \alpha t \nabla f(x)^T \Delta x$$

Hence backtracking line search will always terminate.

[Why $\alpha < \frac{1}{2}$ will be clear in convergence analysis later].

Note: If $\text{dom} f \neq \mathbb{R}^n$ then care must be taken to ensure that $x + t\Delta x \in \text{dom} f$ before checking the inequality in backtracking line search.

3) Gradient descent method:

The question here is how to choose the (descent) search direction Δx ?

A candidate is the negative gradient:

$$\Delta x = -\nabla f(x)$$

Gradient descent method:

given a starting point $x \in \text{dom} f$.

repeat

1. $\Delta x := -\nabla f(x)$

2. if $\|\nabla f(x)\|_2 \leq \epsilon$, break
else

3. Line search: choose t via exact or backtracking

4. Update $x := x + t\Delta x = x - t\nabla f(x)$

end.

t) Convergence analysis: We will show the convergence for backtracking line search type

Assume strong convexity holds: $\exists m, M > 0$ s.t.

$$mI \preceq \nabla^2 f(x) \preceq MI \quad \forall x \in S$$

Consider step length t such that $x - t\nabla f(x) \in S$.

Recall $f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{M}{2} \|y-x\|_2^2 \quad \forall x, y \in S$

Let $y = x - t \nabla f(x)$ then

$$f(x - t \nabla f(x)) \leq f(x) - t \|\nabla f(x)\|_2^2 + \frac{M}{2} t^2 \|\nabla f(x)\|_2^2$$

Now use the inequality

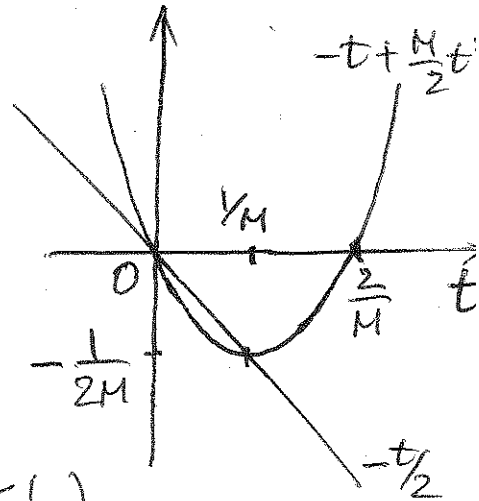
$$-t + \frac{Mt^2}{2} \leq -\frac{t}{2} \quad \text{for } 0 \leq t \leq \frac{1}{M}$$

we have

$$f(x - t \nabla f(x)) \leq f(x) - \frac{t}{2} \|\nabla f(x)\|_2^2$$

$$\leq f(x) - \alpha t \|\nabla f(x)\|_2^2$$

for $0 < \alpha < \frac{1}{2}$. (Hence the range $\alpha < \frac{1}{2}$)



Thus for t small enough ($0 < t \leq \frac{1}{M}$), the backtracking line search will terminate.

This happens either with $t=1$ or for some $1 > t \geq \frac{\beta}{M}$.

Thus we can write

$$f(x^+) \leq f(x) - \min\left\{\alpha, \frac{\beta\alpha}{M}\right\} \|\nabla f(x)\|_2^2$$

$$\rightarrow f(x^+) - p^* \leq f(x) - p^* - \min\left\{\alpha, \frac{\beta\alpha}{M}\right\} \|\nabla f(x)\|_2^2$$

Recall also that

$$f(x) - p^* \leq \frac{1}{2m} \|\nabla f(x)\|_2^2 \quad \forall x \in S.$$

thus

$$f(x^+) - p^* \leq \left(1 - 2m \cdot \min\left\{\alpha, \frac{\beta\alpha}{M}\right\}\right) (f(x) - p^*).$$

Therefore $f(x^{(k)}) - p^* \leq c^k (f(x^{(0)}) - p^*)$

for a constant $c = 1 - \min\left\{2m\alpha, \frac{2\beta\alpha m}{M}\right\} < 1$.

Lecture 18:

+) Example: A quadratic problem

$$f(x) = \frac{1}{2}(x_1^2 + \gamma x_2^2) \quad \gamma > 0.$$

This problem has an analytical solution of $x_1 = x_2 = 0$.

$$\nabla^2 f(x) = \begin{bmatrix} 1 & 0 \\ 0 & \gamma \end{bmatrix} \rightarrow \begin{aligned} m &= \min(1, \gamma) \\ M &= \max(1, \gamma) \end{aligned}$$

Apply gradient search with exact line search, we can obtain closed form expression for the iterates

$$x_1^{(k)} = \gamma \left(\frac{\gamma-1}{\gamma+1} \right)^k, \quad x_2^{(k)} = \left(-\frac{\gamma-1}{\gamma+1} \right)^k$$

$$\text{and } f(x^{(k)}) = \left(\frac{\gamma-1}{\gamma+1} \right)^{2k} f(x^{(0)}) \text{ for } x^{(0)} = (\gamma, 1).$$

Convergence is linear (exactly): the error reduces by the factor $[(\gamma-1)/(\gamma+1)]^2$ at each iteration.

If $\gamma = 1 \rightarrow$ algorithm takes exactly 1 iteration.

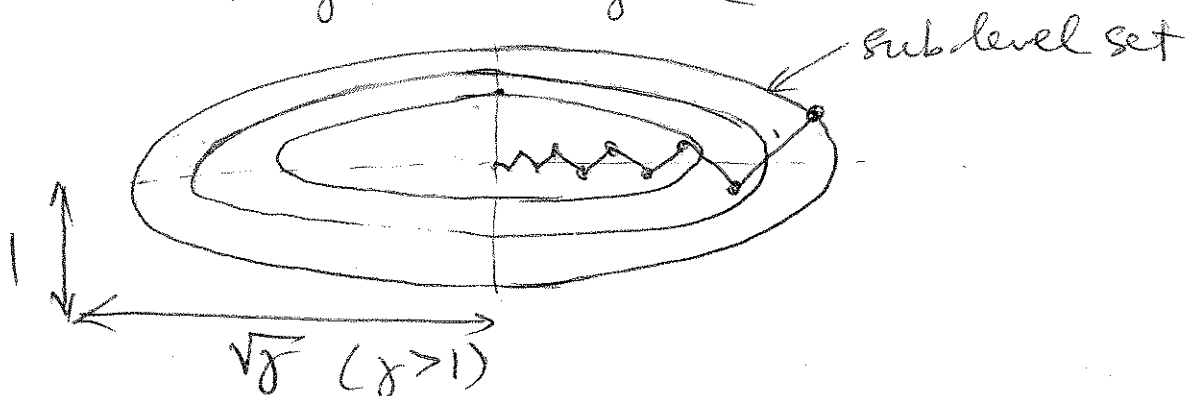
γ around 1 (say $\frac{1}{3} \leq \gamma \leq 3$) \rightarrow convergence is fast.

$\gamma \gg 1$ or $\gamma \ll 1 \rightarrow$ convergence is very slow.

Using the bound derived earlier, error is reduced each iteration by exactly

$$\left(\frac{1 - m/M}{1 + m/M} \right)^2 = \frac{K-1}{K+1}$$

If $K = \frac{M}{m}$ is large \rightarrow convergence is slow.



4) Several observations:

- The gradient method exhibits approximately linear convergence. (error reduces by a factor of c each iteration)
- The choice of backtracking parameters α, β has noticeable but not dramatic effect on the convergence.
- The convergence rate depends strongly on the condition number of the Hessian. Convergence is slow even for $k \approx 100$. For $k \gg 1000$, the gradient method is practically useless.

4. Steepest descent method:

This method determines the "steepest" direction according to some norm.

- First-order Taylor approximation of $f(x+u)$ around x :
$$f(x+u) \approx f(x) + \nabla f(x)^T u$$

$\nabla f(x)^T u$: directional derivatives of f at x in direction u .

- Normalized steepest descent direction:

$$\Delta x_{nsd} = \underset{u}{\operatorname{argmin}} \{ \nabla f(x)^T u \mid \|u\| \leq 1 \}$$

for some norm $\|u\|$.

Δx_{nsd} is the direction in the unit ball in norm $\|u\|$ that extends farthest in the direction $-\nabla f(x)$.

For the 2-norm $\|u\|_2$, steepest descent direction is the same as gradient descent.

- Unnormalized steepest descent direction:

$$\Delta x_{sd} = \|\nabla f(x)\|_* \Delta x_{nsd}$$

where $\|\cdot\|_*$ denotes the dual norm:

$$\nabla f(x)^T \Delta x_{sd} = \|\nabla f(x)\|_* \nabla f(x)^T \Delta x_{ned} = -\|\nabla f(x)\|_*^2$$

Note: A dual norm is given by:

$$\|z\|_* = \sup \{ z^T x \mid \|x\| \leq 1 \}$$

For example: • The dual of the Euclidean norm is the Euclidean norm itself.

• The dual of the l_∞ -norm is the l_1 -norm.

Steepest descent alg.

given a starting point $x \in \text{dom} f$

repeat

1. Compute the steepest descent direction Δx_{sd} .
2. Line search: choose step t
3. Update: $x := x + t \Delta x_{sd}$.

until stopping criterion is satisfied.

Convergence properties are similar to gradient descent (linear convergence)

Examples:

(i) Euclidean norm: $\Delta x_{sd} = -\nabla f(x)$

(ii) Quadratic norm:

$$\|z\|_P = (z^T P z)^{1/2} = \|P^{1/2} z\|_2, \quad P \in S_{++}^n$$

Steepest descent direction:

$$\Delta x_{sd} = -P^{-1} \nabla f(x)$$

This quadratic norm is equivalent to a change of coordinates:

Let $\tilde{u} = P^{1/2} u$ and $\tilde{f}(\tilde{u}) = f(P^{-1/2} \tilde{u}) = f(u)$

then $\nabla \tilde{f}(\tilde{x}) = -P^{-1/2} \nabla f(P^{-1/2} \tilde{x}) = -P^{-1/2} \nabla f(x)$

The gradient search direction corresponds to

$$\Delta x = \bar{P}^{-1/2} (-\bar{P}^{-1/2} \nabla f(x)) = -\bar{P}^{-1} \nabla f(x)$$

for the original variable x .

Thus steepest direction is the same as the ^(negative) gradient direction, after the change of coordinate $\bar{x} = \bar{P}^{1/2} x$.

+) The choice of P can have a strong effect on the rate of convergence.

◦ If the change of coordinate by P reduces the condition number of the resulting sublevel sets \rightarrow speed up convergence, and vice versa.

See figures 9.11 - 9.15 in the text for examples.

5. Newton's method:

◦ For gradient and steepest descent methods, we approximate the objective function at the current value $x^{(k)}$ by a linear function going through $x^{(k)}$.

◦ For Newton's method, we approximate it by a quadratic function through $x^{(k)}$.

+) Newton step:

$$\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$$

For convex functions, $\nabla^2 f(x) \succeq 0$, thus

$$\nabla f(x)^T \Delta x_{nt} = -\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) < 0$$

unless $\nabla f(x) = 0 \rightarrow \Delta x_{nt}$ is a descent direction.

Interpretations:

+) Minimizer of second-order Taylor approximation:

$$\hat{f}(x+u) = f(x) + \nabla f(x)^T u + \frac{1}{2} u^T \nabla^2 f(x) u.$$

$\hat{f}(x+u)$ is a convex quadratic approximation of $f(x)$ at x .

The minimizer of $\hat{f}(x+u)$ is Δx_{nt} :

$$\frac{\partial \hat{f}}{\partial u} = u^T \nabla^2 f(x) + \nabla f(x)^T = 0$$

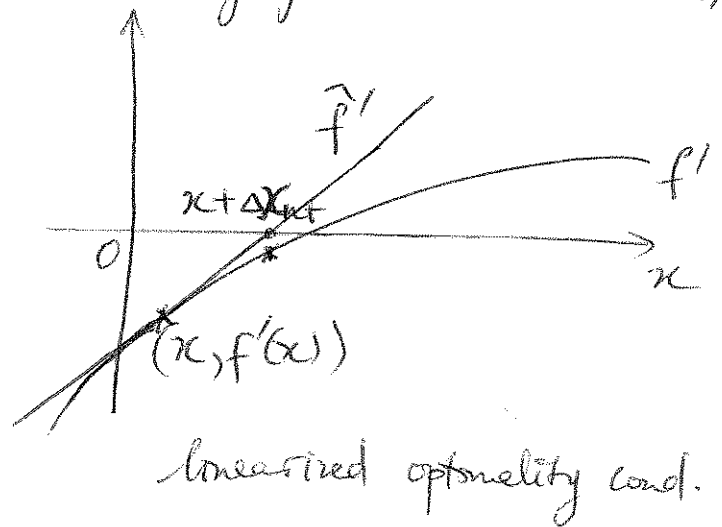
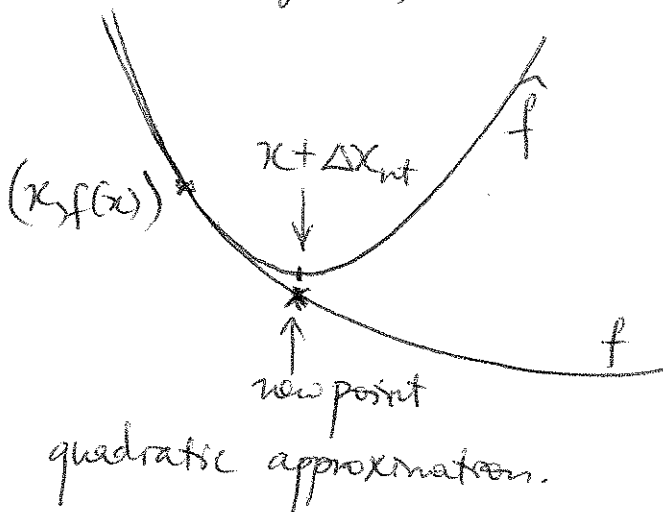
$$\rightarrow u^* = -\nabla^2 f(x)^{-1} \nabla f(x) = \Delta x_{\text{nt}}$$

Since $f(x)$ is twice differentiable, the quadratic model is a very good approximation of $f(x)$ when x is near x^* .

+) Steepest descent direction in Hessian norm:

$$\|u\|_{\nabla^2 f(x)} = (u^T \nabla^2 f(x) u)^{1/2}$$

When x is near x^* , we have $\nabla^2 f(x) \approx \nabla^2 f(x^*)$, which makes the Hessian after this change of coordinate to have almost a condition number of 1 \rightarrow very good search direction (fast convergence)



+) Linearized optimality condition: $\nabla f(x^*) = 0$, linearize $\nabla f(x)$ to get $\nabla f(x+u) \approx \nabla f(x) + \nabla^2 f(x) u = 0 \rightarrow u = \Delta x_{\text{nt}}$

Lecture 19:

+) Affine invariance of the Newton step (property):

If we perform a change of coordinate: $x = Ty$, $T \in \mathbb{R}^n$
 T non-singular.

$$\text{Denote } \tilde{f}(y) = f(Ty) = f(x)$$

$$\rightarrow \nabla \tilde{f}(y) = T^T \nabla f(Ty)$$

$$\nabla^2 \tilde{f}(y) = T^T \nabla^2 f(Ty) T$$

Then the Newton step for $\tilde{f}(y)$ at y is:

$$\Delta y_{\text{nt}} = - (T^T \nabla^2 f(Ty) T)^{-1} T^T \nabla f(Ty)$$

$$= - T^{-1} \nabla^2 f(x)^{-1} \nabla f(x)$$

$$= T^{-1} \Delta x_{\text{nt}}$$

Thus the Newton updates of \tilde{f} and f are also related by the same affine transformation:

$$x + \Delta x_{\text{nt}} = T(y + \Delta y_{\text{nt}})$$

+) Newton decrement:

$$\lambda(x) = (\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x))^{1/2}$$

is called the Newton decrement at x .

• It tells approximately how far point x is from the optimal.

$$f(x) - \inf_y \tilde{f}(y) = f(x) - \tilde{f}(x + \Delta x_{\text{nt}}) = \frac{1}{2} \lambda(x)^2$$

→ $\frac{1}{2} \lambda(x)^2$ gives an estimate of $f(x) - p^*$ (not a lower or upper bound, just an estimate - more later).

• It gives the norm of the Newton step:

$$\lambda(x)^2 = \Delta x_{\text{nt}}^T \nabla^2 f(x) \Delta x_{\text{nt}}$$

Also shows up in backtracking linesearch:

$$\nabla f(x)^T \Delta x_{\text{nt}} = -\lambda(x)^2$$

$\lambda(x)$ is affine invariant

$$\lambda(y) = \lambda(x) \text{ for } \tilde{f}(y) = f(Ty) = f(x) \\ x = Ty$$

→ Newton's method:

given a starting point $x \in \text{dom} f$, tolerance $\varepsilon > 0$
repeat

1. Compute the Newton step and decrement

$$\Delta x_{\text{nt}} := -\nabla^2 f(x)^{-1} \nabla f(x)$$

$$\lambda^2 := \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x)$$

2. Check stopping criterion: quit if $\lambda^2 \leq \varepsilon$.

3. Linesearch: Choose step size t by backtracking linesearch.

4. Update: $x := x + t \Delta x_{\text{nt}}$.

→ Convergence analysis:

• Again assume strong convexity: $\exists m, M > 0$ s.t.

$$mI \preceq \nabla^2 f(x) \preceq MI \quad \forall x \in S.$$

• Also assume the Hessian of f is Lipschitz continuous on S with constant $L > 0$:

$$\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L \|x - y\|_2 \quad \forall x, y \in S.$$

Intuitively, L represents a bound on the third derivative of f .

$L = 0$ for a quadratic function

L small \rightarrow quadratic approximation of f at x is good

L large \rightarrow quadratic approximation is poor.

→ L will play an important role in the convergence rate.

+) Idea of convergence proof:

Show $\exists \eta$ and γ :

$$0 < \eta \leq \frac{m^2}{L}, \quad \gamma > 0$$

such that there are 2 regions of convergence:

• If $\|\nabla f(x^{(k)})\|_2 \geq \eta$ then

$$f(x^{(k+1)}) - f(x^{(k)}) \leq -\gamma. \quad (\text{damped Newton phase}).$$

• If $\|\nabla f(x^{(k)})\|_2 < \eta$, then

$$\frac{L}{2m^2} \|\nabla f(x^{(k+1)})\|_2 \leq \left(\frac{L}{2m^2} \|\nabla f(x^{(k)})\|_2 \right)^2. \quad \text{quadratic convergence}$$

+) Damped Newton phase:

- Function value decreases by at least γ at each iteration.
- Line search with backtracking is used
- Number of iterations in this step:

$$\# \text{ iterations} \leq \frac{f(x^{(0)}) - f^*}{\gamma}$$

+) Quadratic convergence phase:

• Since $\eta \leq \frac{m^2}{L}$, once the algorithm enters this phase, it will stay in this phase:

$\forall l > k: \|\nabla f(x^{(l)})\|_2 < \eta$ since

$$\begin{aligned} \|\nabla f(x^{(k+1)})\|_2 &\leq \frac{2m^2}{L} \cdot \left(\frac{1}{2} \cdot \frac{L}{m^2} \|\nabla f(x^{(k)})\|_2 \right)^2 \\ &\leq \left(\frac{1}{2} \cdot \frac{1}{\eta} \right) \cdot \eta^2 = \frac{1}{2} \eta < \eta \end{aligned}$$

• Apply the inequality recursively, we get:

$$\frac{L}{2m^2} \|\nabla f(x^{(k)})\|_2 \leq \left(\frac{L}{2m^2} \|\nabla f(x^{(k)})\|_2 \right)^2 \leq \left(\frac{1}{2} \right)^{2^{l-k}}$$

(since $\|\nabla f(x^{(k)})\|_2 < \eta \leq \frac{m^2}{L}$)

Thus

$$f(x^{(l)}) - p^* \leq \frac{1}{2m} \|\nabla f(x^{(l)})\|_2 \leq \frac{2m^3}{L^2} \left(\frac{1}{2} \right)^{2^{l-k+1}}$$

↑ Strong convexity
 ↓ ϵ_0

→ Convergence is extremely rapid in this phase.

+) The overall number of iterations is bounded above by

$$\frac{f(x^{(0)}) - p^*}{\gamma} + \log_2 \log_2 \left(\frac{\epsilon_0}{\epsilon} \right)$$

$$\epsilon_0 = \frac{2m^3}{L}$$

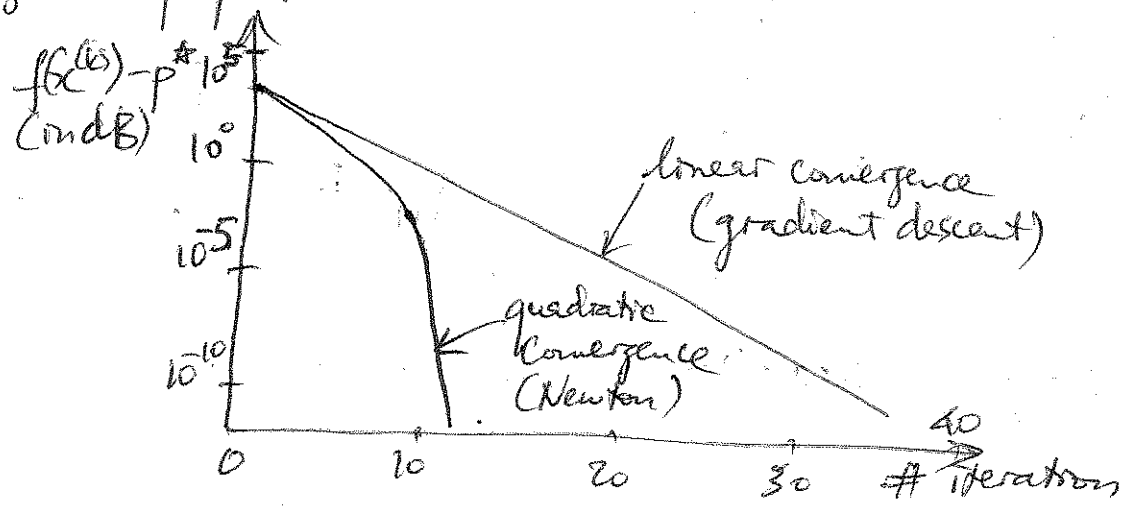
ϵ : prespecified tolerance level.

For example, six iterations in the quadratic phase gives an accuracy of about

$$\epsilon \approx 5 \times 10^{-20} \epsilon_0$$

→ the quadratic phase can be considered to have almost constant # of iterations (≈ 6).

• In practice, the constants m, L (hence γ, ϵ_0) are usually unknown. But this convergence analysis provides insight into convergence properties.



+) More detailed convergence proof:

o Damped Newton phase:

By strong convexity $\nabla^2 f(x) \preceq MI \quad \forall x \in S$:

$$\begin{aligned} f(x + t \Delta x_{\text{nt}}) &\leq f(x) + t \nabla f(x)^T \Delta x_{\text{nt}} + \frac{M}{2} \|\Delta x_{\text{nt}}\|_2^2 t^2 \\ &\leq f(x) - t \lambda(x)^2 + \frac{M}{2m} t^2 \lambda(x)^2 \left(f(x) - \frac{t}{2} \lambda \right) \end{aligned}$$

Noting step size $\hat{t} = \frac{\eta}{M}$ satisfies backtracking line search, then line search must return $t \geq \beta \frac{\eta}{M}$, thus

$$\begin{aligned} f(x^+) - f(x) &\leq -\frac{t}{2} \lambda(x)^2 \leq -\alpha t \lambda(x)^2 \\ &\leq -\frac{\alpha \beta \eta}{M^2} \|\nabla f(x)\|_2^2 \leq -\alpha \beta \frac{\eta}{M^2} \eta^2 \end{aligned}$$

[Here we use $\lambda(x)^2 = \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \geq \frac{1}{M} \|\nabla f(x)\|_2^2$

Thus choose $\gamma = \alpha \beta \frac{\eta}{M^2} \eta^2$ satisfies the damped phase.

o Quadratic phase:

For this phase we can always use step size $t=1$ in backtracking search. We skip the proof for this here.

Now $x^+ = x + \Delta x_{\text{nt}}$ and

$$\begin{aligned} \|\nabla f(x^+)\|_2^2 &= \|\nabla f(x + \Delta x_{\text{nt}}) - \nabla f(x) - \nabla^2 f(x) \Delta x_{\text{nt}}\|_2 \\ &= \left\| \int_0^1 \nabla^2 f(x + t \Delta x_{\text{nt}}) \cdot \Delta x_{\text{nt}} dt - \nabla^2 f(x) \Delta x_{\text{nt}} \right\|_2 \\ &= \left\| \int_0^1 [\nabla^2 f(x + t \Delta x_{\text{nt}}) - \nabla^2 f(x)] \Delta x_{\text{nt}} dt \right\|_2 \end{aligned}$$

$$\begin{aligned} &\stackrel{\text{Lipschitz}}{\leq} \frac{L}{2} \|\Delta x_{\text{nt}}\|_2^2 = \frac{L}{2} \|\nabla^2 f(x)^{-1} \nabla f(x)\|_2^2 \\ &\stackrel{\text{Strong convexity}}{\leq} \frac{L}{2m^2} \|\nabla f(x)\|_2^2 \end{aligned}$$

Lecture 20:

+) More detailed convergence proof:

o Damped Newton phase:

By strong convexity $\nabla^2 f(x) \preceq M I \quad \forall x \in S$:

$$\begin{aligned} f(x + t \Delta x_{nt}) &\leq f(x) + t \nabla f(x)^T \Delta x_{nt} + \frac{M}{2} \|\Delta x_{nt}\|_2^2 t^2 \\ &\leq f(x) - t \lambda(x)^2 + \frac{M}{2m} t^2 \lambda(x)^2 \left(f(x) - \frac{t}{2} \lambda t \right) \end{aligned}$$

Noting step size $\hat{t} = \frac{\eta}{M}$ satisfies backtracking line search, then line search must return $t \geq \beta \frac{\eta}{M}$, thus

$$\begin{aligned} f(x^+) - f(x) &\leq -\frac{t}{2} \lambda(x)^2 \leq -\alpha t \lambda(x)^2 \\ &\leq -\frac{\alpha \beta m}{M^2} \|\nabla f(x)\|_2^2 \leq -\alpha \beta \frac{m}{M^2} \eta^2 \end{aligned}$$

[Here we use $\lambda(x)^2 = \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \geq \frac{1}{M} \|\nabla f(x)\|_2^2$

Thus choose $\gamma = \alpha \beta \frac{m}{M^2} \eta^2$ satisfies the damped phase.

o Quadratic phase:

For this phase we can always use step size $t=1$ in backtracking search. We skip the proof for this here.

Now $x^+ = x + \Delta x_{nt}$ and

$$\begin{aligned} \|\nabla f(x^+)\|_2^2 &= \|\nabla f(x + \Delta x_{nt}) - \nabla f(x) - \nabla^2 f(x) \Delta x_{nt}\|_2 \\ &= \left\| \int_0^1 \nabla^2 f(x + t \Delta x_{nt}) \cdot \Delta x_{nt} dt - \nabla^2 f(x) \Delta x_{nt} \right\|_2 \\ &= \left\| \int_0^1 [\nabla^2 f(x + t \Delta x_{nt}) - \nabla^2 f(x)] \Delta x_{nt} dt \right\|_2 \end{aligned}$$

$$\begin{aligned} &\stackrel{\text{Lipschitz}}{\leq} \frac{L}{2} \|\Delta x_{nt}\|_2^2 = \frac{L}{2} \|\nabla^2 f(x)^{-1} \nabla f(x)\|_2^2 \\ &\stackrel{\text{Strong convexity}}{\leq} \frac{L}{2m^2} \|\nabla f(x)\|_2^2 \end{aligned}$$

Can show $\eta = \min \left\{ 1, \frac{3(1-2\alpha)}{L} \right\} \frac{m^2}{L}$.

Summary for Newton's method:

- Convergence rapid, quadratic near x^* .
- Affine invariant, thus insensitive to choice of coordinate or condition number of sublevel sets.
- The condition number of sublevel sets (or the Hessian) only affects the numerical inversion of the Hessian, but has little effect on the rate of convergence. Newton's method can tolerate large condition numbers of up to 10^{10} , whereas gradient descent tolerates far smaller numbers (practically useless for $K > 1000$).
- Newton's method scales well with problem size. For example, for problems in \mathbb{R}^{10} and \mathbb{R}^{1000} the number of iterations can be comparable.

⊖ The difficulty of the Newton's method is the cost of forming (computing) and storing the Hessian, and the cost of computing the Newton step, requires to solve $\nabla^2 f(x) \Delta x = -\nabla f(x)$.
Roughly the cost of computing the inverse of an $n \times n$ Hessian matrix is of order n^3 . But for some problem we can exploit the structure of the Hessian to reduce this cost.

+ Some variants of Newton's method:

The goal is to reduce the computational complexity of Newton's.

• Quasi-Newton methods: Replace $\nabla^2 f(x)$ by approximation H . Many update rules for H , all satisfies:

$$\bullet H = H^T \succ 0$$

$$\bullet \text{secant condition: } \nabla f(x^+) - \nabla f(x) = H^+(x^+ - x)$$

$$\bullet H^+ \nabla f(x) \text{ is more easily computed than } \nabla^2 f(x)^{-1} \nabla f(x)$$

Broyden-Fletcher-Goldfarb-Shanno (BFGS): (most common)

$$y = \nabla f(x^+) - \nabla f(x) \quad ; \quad s = x^+ - x$$

$$\text{then } H^+ = H + \frac{yy^T}{y^T s} - \frac{Hs s^T H}{s^T H s} \rightarrow O(n^2)$$

→ Self-concordance:

• Motivation:

- Newton's method is affine invariant, but the convergence analysis is not.

- Often do not know constants m, M, L in practice

- Constants m, M, L can depend on starting point.

• Self-concordance condition (Nesterov & Nemirovski).

allows a new analysis of Newton's method

- is affine invariant

- involves no unknown constants

- is valid for many functions f , including the logarithmic barrier function (more on this later in interior-point methods)

• Self-concordance condition:

convex $f: \mathbb{R} \rightarrow \mathbb{R}$ is self-concordance if

$$|f'''(x)| \leq 2f''(x)^{3/2} \quad \forall x \in \text{dom} f.$$

Examples:

• $f(x) = -\log x$ are SC.

• $f(x) = x \log x - \log x$

• $f(x) = -\log \det x$

• Some simple properties:

- affine invariant:

$$f(x) \text{ SC} \Leftrightarrow g(z) = f(Az + b) \text{ is SC.}$$

- sum and scaling:

$$f, \tilde{f} \text{ SC} \rightarrow f + \tilde{f} \text{ SC}$$

$$f \text{ SC} \rightarrow \alpha f \text{ SC.}$$

Thus: $-\sum_i \log(b_i - a_i^T x)$ is SC

$-\log \det(F_0 + x_1 F_1 + \dots + x_n F_n)$ is SC.

+ Analysis of Newton's method for SC functions:

Can show that with backtracking or exact line search:

$$\# \text{ iterations} \leq \frac{f(x^{(0)}) - f^*}{\eta_2} + \log_2 \log_2 \left(\frac{2}{\epsilon} \right)$$

where η_2 depends only on backtracking line search parameters:

$$\eta_2 = \beta \frac{\alpha(1/2 - \alpha)^2}{5 - 2\alpha}$$

+ SC functions allow a more explicit analysis of convergence and complexity.

It is not known if SC functions are easier to minimize than non-SC functions.